



# Módulo I

# Fundamentos Recuperación en Internet

OpenCourseWare

Recuperación y Acceso a la Información

# Contenidos

- Buscadores Web
- Tipología
- Internet invisible
- Otros tipos de buscadores
- Tendencias en Internet

# ¿Qué es un Buscador?

Un buscador es un software que busca en una base de datos o repositorio documental, conforme a algunos criterios y prioridades específicas.

## Objetivos

- Indizar la red constantemente para permitir la consulta de sus recursos
- Encontrar los documentos que contengan las palabras clave introducidas por el usuario

## Características Web

- Volumen
- Obsolescencia
- Actualización
- Tipología
- Tráfico

## Tipos

- Directorios
- Motores de búsqueda
- Meta-buscadores

# Contenidos

- Buscadores Web
- Tipología
- Internet invisible
- Otros tipos de buscadores
- Tendencias en Internet

# Directorios

¿Qué son?

– Páginas con clasificación intelectual por temas

• Características

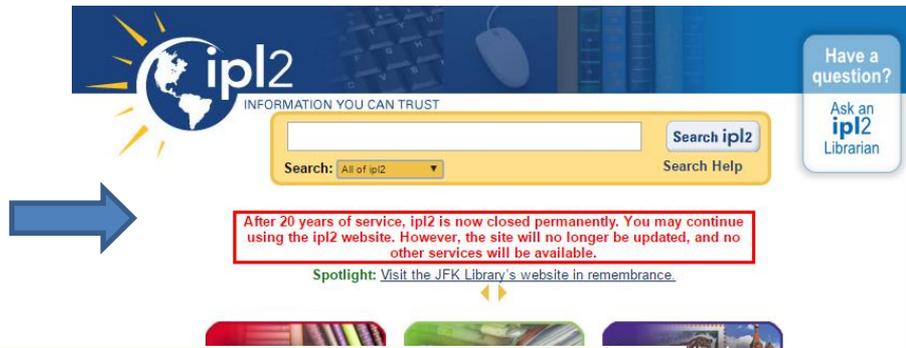
- Actualización
- Exhaustividad
- Relevancia
- Calidad

– Ejemplos: DMOZ (1998-2017), Yahoo (1994-2014), Ipl2 (1995-2013), Categorización Wikipedia...



Directorio Yahoo!

# Directorios

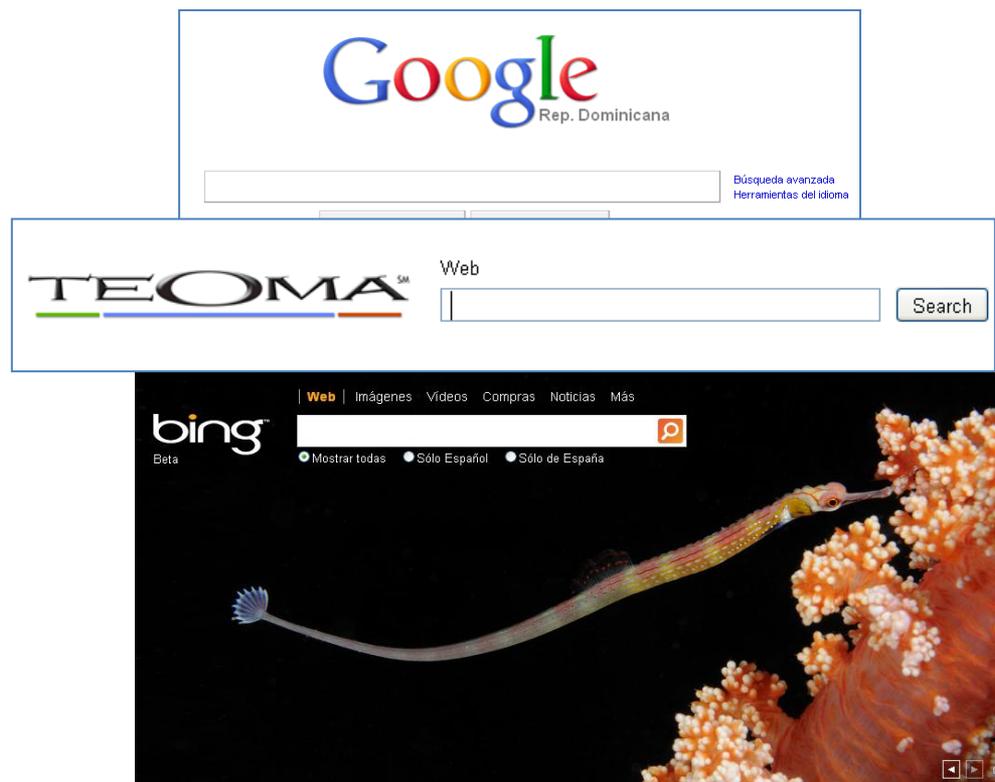


# Motores de búsqueda

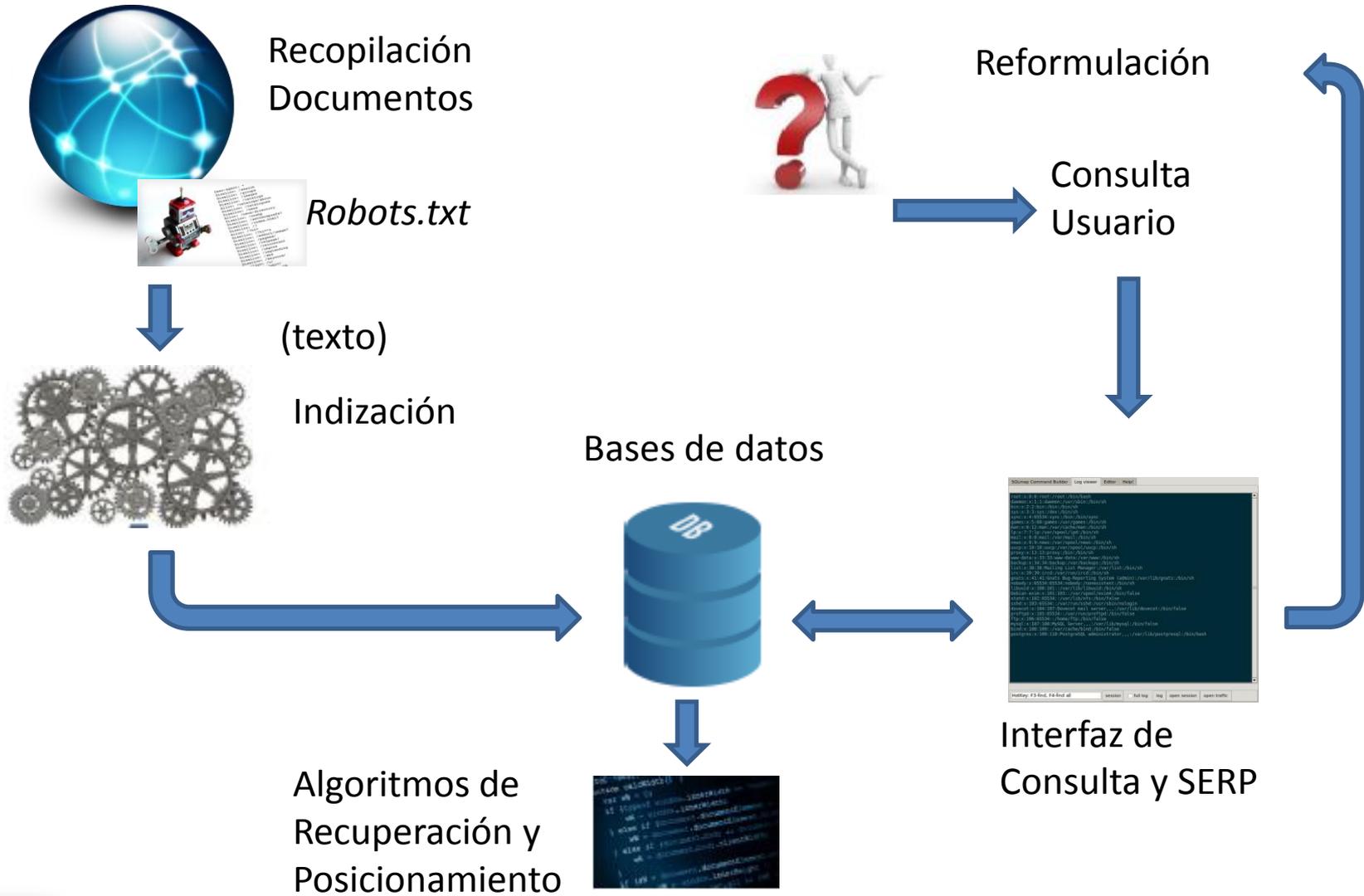
- ¿Qué son?
  - Aplicación de recolección y búsqueda

- Características
  - Exhaustividad
  - Actualización
  - Manipulación
  - Problemas léxicos
  - Seguridad

- Ejemplos: Google



# Sistema de Recuperación



SERP (Search Engine Result Page)

# Estadísticas de uso motores

Search Engine	Desktop	Smart phone
Google	77.82%	93.76%
Bing	7.96%	1.06%
Yahoo! Search	6.39%	4.01%
Baidu	6.37%	0.37%
Ask	0.15%	0.04%
AOL	0.08%	0.01%

Fuente <https://www.netmarketshare.com/>, Dic. 2016

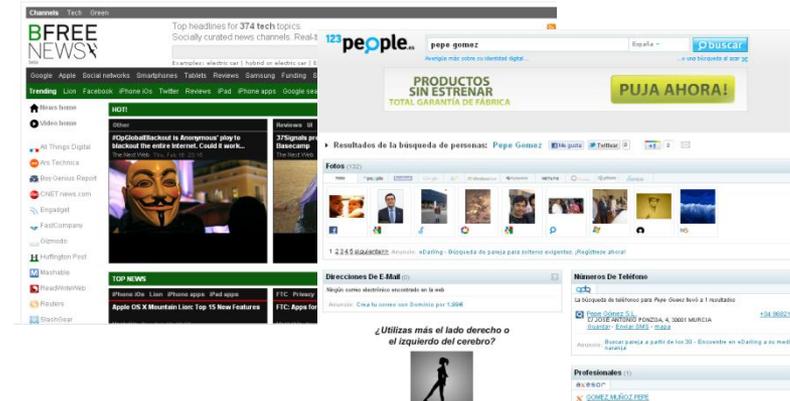
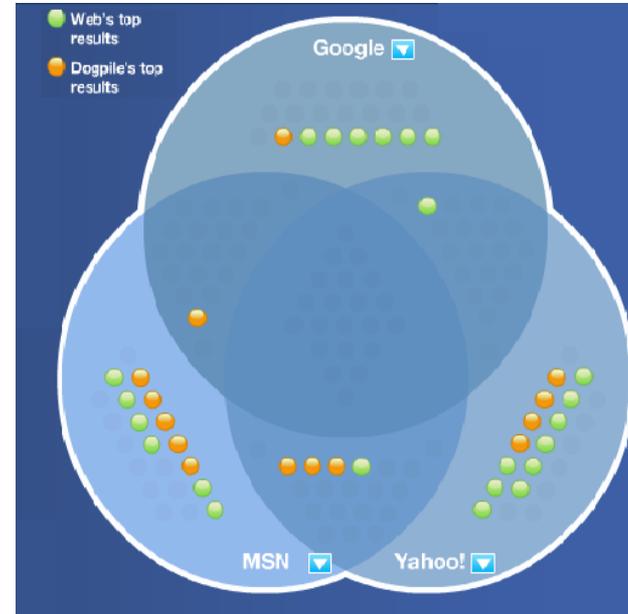
Hay diferencias regionales:

Baidu (China, 55% y Qihoo 360, 28%),  
 Yandex (Rusia, 58%),  
 Corea del Sur (Naver, 77%),  
 Japon (Yahoo Japan, 40%, tras Google)

Fuente <http://returnnonnow.com/> 2015

# Metabuscadores

- ¿Qué son?
  - Aplicación que agrega bases de datos, motores y/o directorios
- Características
  - **Algoritmo**
  - **Políticas de uso**
- Diferencia entre:
  1. Mashups
  2. Multibuscadores, agentes de búsqueda, metabuscadores



Ejemplos: Rastreator, Rumbo, Kayak, ...

# Información almacenada en Google

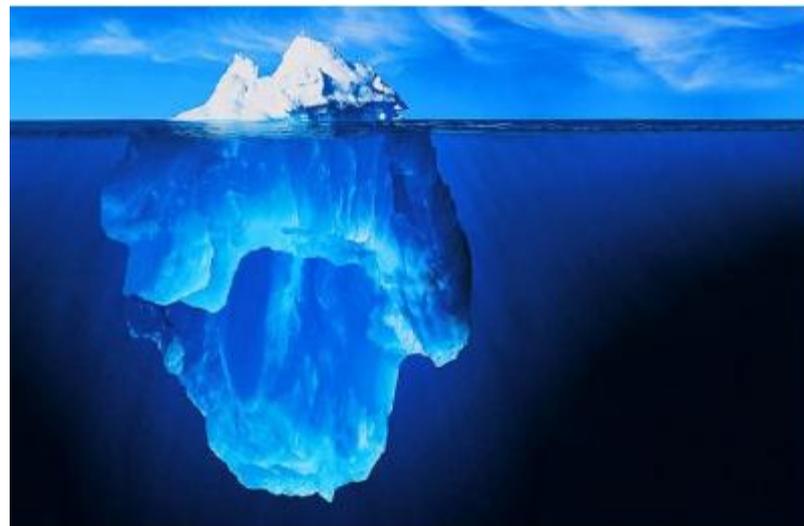
- **Site**
- **Allinurl**
- **Allintitle**
- **Link (no fiable)**
- **Info**
- **|**
- **“”**
- **-**
- **111..333**
- **Define**
- **Cache**
- **Press**
- **Location**
- **Index of**
- **Filetype**
- **Related**

# Contenidos

- Buscadores Web
- Tipología
- Internet invisible
- Otros tipos de buscadores
- Tendencias en Internet

# Problemas de los buscadores: Internet Invisible

- Carencias de motores de búsqueda debido a:
  - Formatos no soportados
  - Necesidad de validación
  - Formularios de acceso
  - Sitos excluidos expresamente
  - Páginas creadas automáticamente
  - No textual
- No confundir con la DeepWeb
- Aproximadamente tiene un 50% más de tráfico que el visible
- Buscadores especializados: p.e. Easy searcher, OPAC, ...



# Contenidos

- Buscadores Web
- Tipología
- Internet invisible
- Otros tipos de buscadores
- Tendencias en Internet

# Recuperación de imágenes

- Texto asociado
- Asignación intelectual
- Reconocimiento
- Ejemplos: Google images, Flickr

**Google**  
Image Labeled BETA **Google Image Labeled**

Welcome to **Google Image Labeled**, a feature of Google Search that allows you to label images and help improve the quality of Google's image search results.

Your nickname: [guest](#) - [Change](#)

[Start labeling](#)

**How does it work?**

You'll be randomly paired with a partner who's online and using the feature. Over a two-minute period, you and your partner will:

- View the same set of images.
- Provide as many labels as possible to describe each image you see.
- Receive points when your label matches your partner's label. The number of points will depend on how specific your label is.
- See more images until time runs out.

After time expires, you can explore the images you've seen and the websites where those images were found. And we'll show you the points you've earned throughout the session.

**Tips:**

**flickr** As You Said

Inicio La vista Crear cuenta Explorar Subir fotos

Explorar

Explora las fotos interesantes de Flickr al elegir un punto en el tiempo.

Seleccionar un mes  
Elegir

Más lugares para explorar:

- Fotos interesantes de los últimos 7 días
- Vista calendario de este mes
- Un mapa del mundo
- Buscador de cámaras
- Cargos más recientes
- Videos en Flickr
- Expos
- App Garden
- El Blog de Flickr

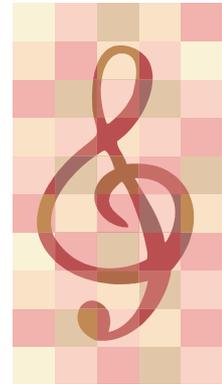
Today's Top Pairs	
1. JBL - guest	1800
2. guest - JBL	1480
3. HarlezGMa - guest	1410
4. Kact - guest	1290
5. guest - guest	1240

All-time Top Contributors	
1. DeS	45663270
2. PS	39999990
3. Zippy	33572910
4. Mw	29476630
5. FrankD	26666660



# Recuperación de música

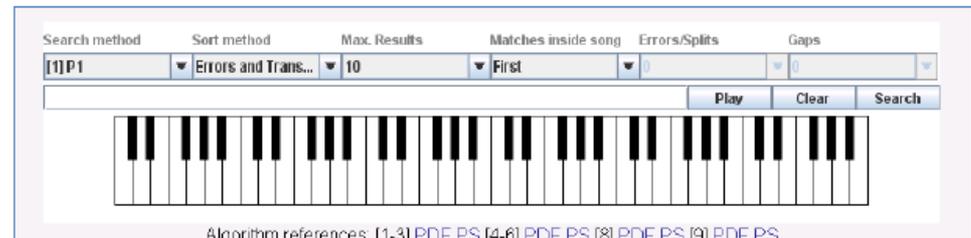
- Similar a texto en origen
- Uso de texto asociado
- Dos tipos: Audio (wav, mp3...) o Simbólica (midi, musicXML...)
- Problemática
  - Especificar consultas
  - Mostrar resultados
  - Derechos



R.Typeke et.al., "A Survey of Music Information Retrieval Systems", ISMIR, 2005

# Recuperación de música (II)

- Tiene muchas aplicaciones:
  - Identificación
  - Detección de plagio
  - Recomendación
  - Generación de pentagramas
- Hay sistemas a escala industrial, pero la mayoría sigue siendo experimental sobre colecciones pequeñas
  - Shazam
  - C-Brahms
  - Liveplasma



# Contenidos

- Buscadores Web
- Tipología
- Internet invisible
- Otros tipos de buscadores
- Tendencias en Internet

# Sistemas pregunta respuesta

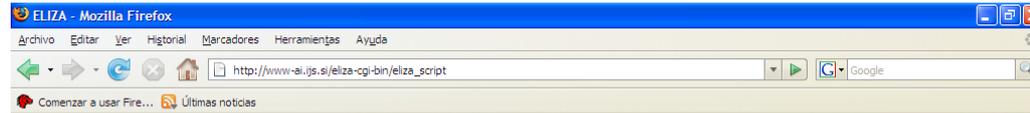
Un sistema pregunta-respuesta devuelve una respuesta concreta a una pregunta y no un conjunto de documentos que puedan contener la respuesta, un subtipo son los asistentes virtuales

- Buscadores con interpretación del lenguaje natural:
  - Asistentes virtuales
  - Answers.com <http://www.answers.com/>
- Buscadores con interpretación semántica
  - Start <http://start.csail.mit.edu/>
  - WolframAlpha <http://www.wolframalpha.com/>

# Asistentes Virtuales

Con:

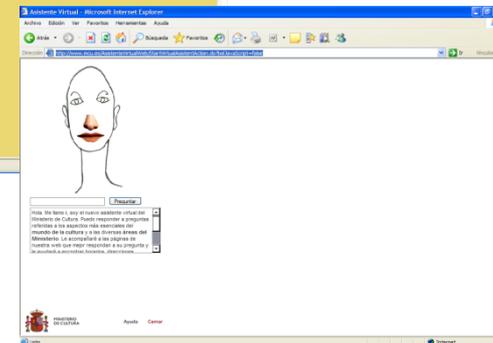
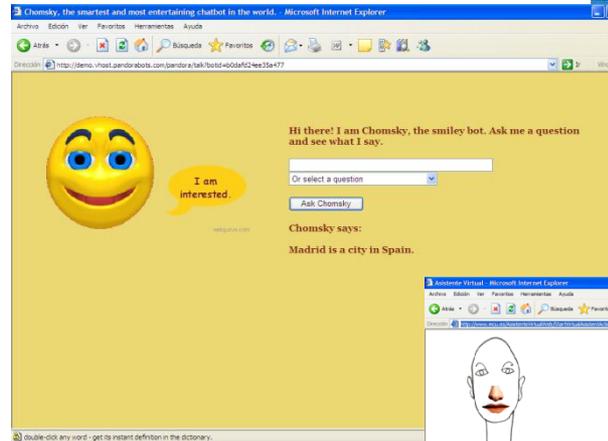
- PLN
- Organización léxica
- Patrones
- Estadísticas



## ELIZA - a friend you could never have before

You: i wish a new car

Eliza: Can you elaborate on that?



- [Ella](#)
- [Chomsky](#)
- [IKEA \(Anna\)](#)
- [Siri \(iOS\)](#)
- [Renfe \(Irene\)](#)

# ¿Qué es la web semántica?



Tim Berners-Lee

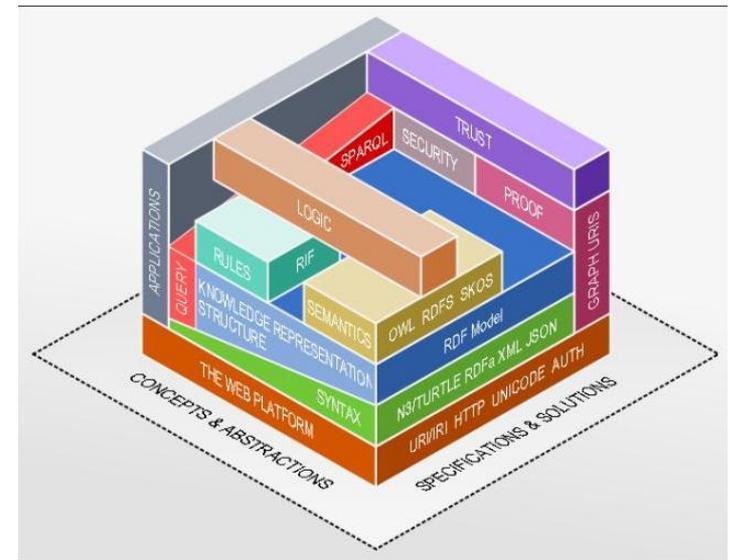
Una Web Semántica es una red de datos que pueden ser procesados directa o indirectamente por máquinas.

Es una web extendida que permitirá a humanos y máquinas trabajar en cooperación mutua.

Guía breve de la Web Semántica <http://www.w3c.es/Divulgacion/Guiasbreves/WebSemantica>

# ¿Para qué sirve la Web Semántica?

Procesar contenido, razonar, combinarlo y realizar deducciones lógicas para resolver problemas cotidianos de forma automática



- Resolución de problemas con ciertas consultas
  - Ambigüedad
  - Variación lingüística
  - Contexto de la consulta
- Problemas
  - Usabilidad
  - Lenguaje natural y codificación semántica
  - Diferentes codificaciones semánticas

# Web Semántica

## Intercambio no ambiguo de información

- Para intercambiar información se necesita Interoperabilidad, compartiendo:
  - Un mismo medio con la misma forma de comunicación
    - Textos con lenguaje de marcado en el Web, p.e. HTML
  - Compartir un **Vocabulario**
    - Vocabularios de Metadatos
  - Compartir una Sintaxis y una misma forma de expresarse
    - p.e. XML/RDF
- Compartir un mismo conocimiento común
  - Tesoros y ontologías



# Contexto: Web Semántica y RDF

La Web Semántica permite navegar a través de datos y semánticas. Requisitos:

1. Identidad
2. Accesibilidad (p.e. protocolos http)
3. Estructura (normalización, RDF)
4. Navegación (recursos enlazados)

*El valor puede ser una URI o un literal*



*RDF está compuesto por tripletas con la forma:*

**SUJETO- ATRIBUTO- VALOR**

## Tripleta

### Sujeto

www.john.com

### Atributo

Autor (URI)

### Valor

John (URI/Literal)

# Diferencia con XML

## Posibilidades con XML:

```
<Libro idTitulo="El_Quijote"
xmlns="http://www.ejemplo.org/libro">
  <autor>Cervantes</autor>
</Libro>
```

```
<obra>
  <libro>El Quijote</libro>
  <autor>Cervantes</autor>
</obra>
```

```
<autor nombre="Cervantes">
  <Libro>El Quijote</Libro>
</autor>
```

## Pero con RDF...

```
<rdf:RDF rdf:ID="El_Quijote"
  xmlns:rdf=http://www.w3.org/1999/02/22-rdf-          syntax-ns#
  xmlns:dc="http://purl.org/dc/elements/1.1/"
  xml:base="http://www.mi_ejemplo.org/libro#">
  <rdf:Description ID="quijote">
    <dc:title>El Quijote</dc:title>
    <dc:creator>Miguel de Cervantes</dc:creator>
  </rdf:Description>
</rdf:RDF>
```

# Namespaces

## Ejemplo con Dublin Core y el elemento Title

```
<?xml version='1.0' encoding='utf-8' standalone='yes'>
```

```
<Libros
```

```
xmlns:dc="http://purl.org/dc/elements/1.1/title">
```

namespace

<http://purl.org/dc/elements/1.1/>

```
<Libro>
```

```
<dc:title lang="esp">Don Quijote de la Mancha</title>
```

```
<Autor>Miguel de Cervantes</Autor>
```

```
</Libro>
```

```
<Libro>
```

```
<dc:title lang="es">La vida es sueño</title>
```

```
<Autor>Calderón de la Barca</Autor>
```

```
</Libro>
```

Dublin Core Metadata Element Set, Version 1.1	
Identifier:	<a href="http://dublincore.org/documents/2008/01/14/dces/">http://dublincore.org/documents/2008/01/14/dces/</a>
Supersedes:	<a href="http://dublincore.org/documents/2006/12/18/dces/">http://dublincore.org/documents/2006/12/18/dces/</a>
Label:	Subject
Definition:	The topic of the resource.
Comment:	Typically, the subject will be represented using keywords, key phrases, or describe the spatial or temporal topic of the resource, use the Coverage element.
Term Name: title	
URI:	<a href="http://purl.org/dc/elements/1.1/title">http://purl.org/dc/elements/1.1/title</a>
Label:	Title
Definition:	A name given to the resource.
Comment:	Typically, a Title will be a name by which the resource is formally known.
Term Name: type	
URI:	<a href="http://purl.org/dc/elements/1.1/type">http://purl.org/dc/elements/1.1/type</a>

Qname o prefijos del namespaces

# Web 2.0

- Web 1.0, no interacción, estática
- Web 2.0
  - Interactividad: Entornos sociales y cooperativos. Espacios de opinión
  - Publicación más simple y estándar. Plantillas, lenguajes, código abierto, ...
  - Descentralizado y flexible
  - Difusión y sindicación mejorada

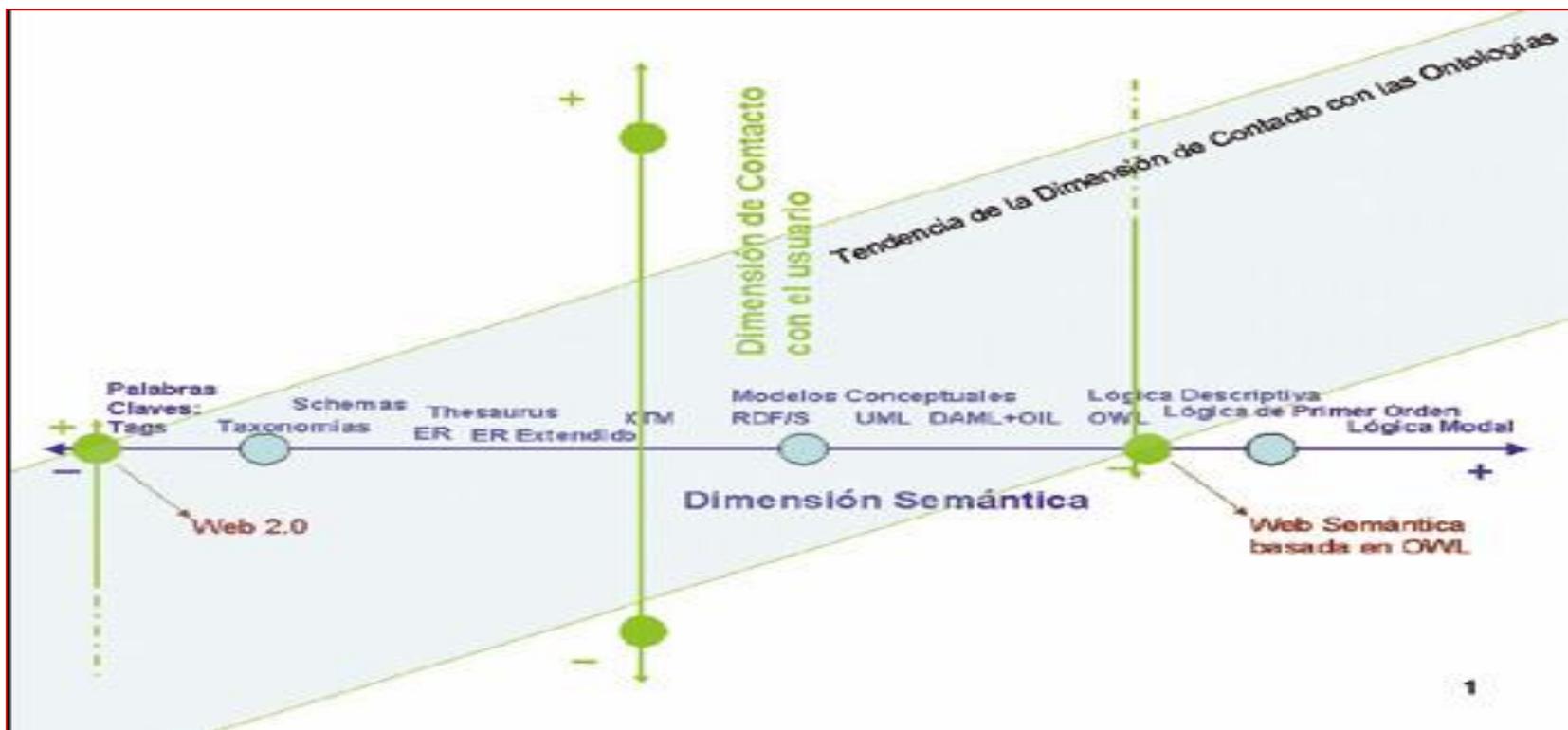


# Web 2.0 vs. Web Semántica

	Web 2.0	Web Semántica
<b>Origen</b>	Constatación de la evolución natural de la Web	Propuesta de Tim Berners para evolucionar la Web
<b>Implantación</b>	Muy alta	Escasa
<b>Coordinación</b>	No existe	Centralizada, sobre todo por el W3C
<b>Foco</b>	<b>Personas</b>	<b>Aplicaciones informáticas</b>
<b>Creación</b>	2003, 1ª conferencia 2004	1999 (Berners-Lee, 1999)

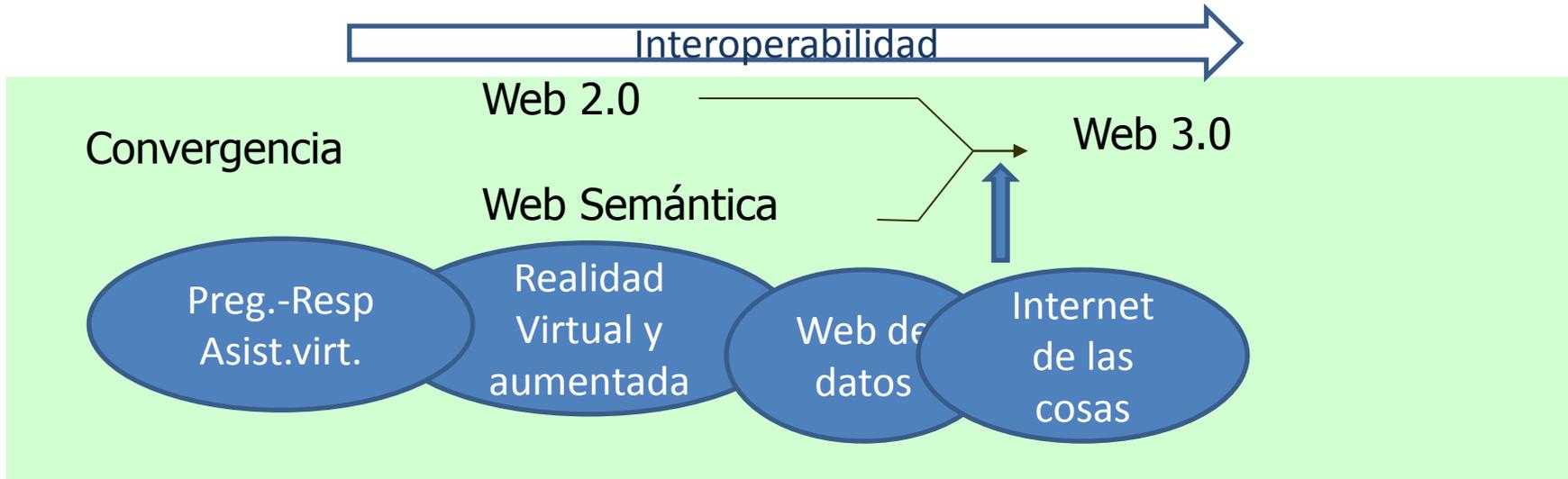
- Entran en confrontación por:
  - Las ontologías son poco legibles (RDF y OWL) por personas y costosas de crear. Las folksonomías son difíciles de interpretar por aplicaciones (polisemia y ambigüedad) pero su creación tiene bajo coste y esfuerzo
  - No hay herramientas de la Web Semántica amigables para los usuarios. Los recursos de la Web Social no son amigables para las aplicaciones
  - Técnicas automáticas de creación de ontologías inmaduras. Problemas de mapeo
  - Duplicidades de vocabularios de metadatos (p.e. SKOS-Core, los PSI, Zthes y MADS)

# Web 2.0 vs. Web Semántica

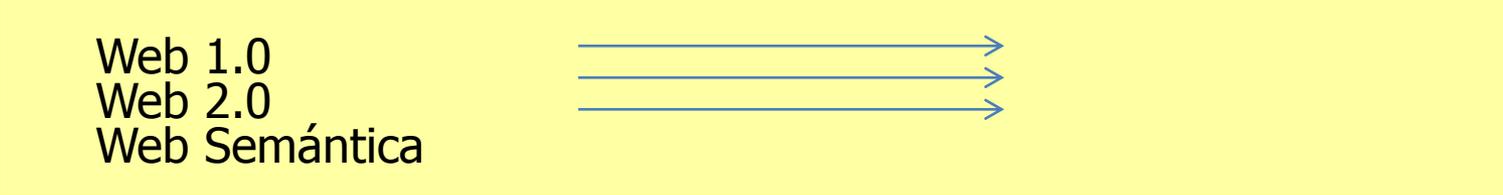


- La Web 3.0 intenta eliminar las dificultades de comunicación con la Web Semántica.
- La web de datos da mayor conocimiento a las aplicaciones web.
- La página web pasa a tener un papel menor, los datos se pueden usar y visualizar de distintas formas.
- El usuario puede leer, publicar pero también ejecutar nuevas aplicaciones.

# Posibilidades evolución Web



La **Web 3.0 es denominada por algunos Web Semántica**, con la necesaria mejora en usabilidad. Los Web sites perderán protagonismo



Siempre debe utilizarse la Web que provee la solución más económica con funcionalidad suficiente

¿Pero y la credibilidad y fiabilidad de los datos?



# Módulo I

## Fundamentos Recuperación en Internet

Colaboradores

J.Morato, V.Palacios

M.Marrero, S.Sánchez-Cuadrado, J.Urbano